

1 **Comparative genomic analysis reveals novel phylogenetically intermediate**  
2 **Streptococci with high phenotypic diversity in the human distal lung microbiota.**

3 Slipa Kanungo<sup>e\*</sup>, Germán Bonilla-Rosso<sup>a\*</sup>, Garance Sarton-Lohéac<sup>a</sup>, Marianne Kuffer<sup>c</sup>,  
4 Markus Hilty<sup>c</sup>, Thomas Geiser<sup>d,e</sup>, Philipp Engel<sup>a</sup>, Sudip Das<sup>a,b,d,e #</sup>

5 <sup>a</sup> Department of Fundamental Microbiology, University of Lausanne, Lausanne, Switzerland

6 <sup>b</sup> Department of Respiratory Medicine, Lausanne University Hospital and University of  
7 Lausanne, Lausanne, Switzerland

8 <sup>c</sup> Institute for Infectious Diseases, University of Bern, Switzerland

9 <sup>d</sup> Department for Pulmonary Medicine, Allergology and Clinical Immunology, Inselspital,  
10 Bern University Hospital, University of Bern, Switzerland

11 <sup>e</sup> Lung Precision Medicine (LPM), Department for BioMedical Research (DBMR), University  
12 of Bern, Switzerland

13

14 Running title: New commensal Streptococci from human lung microbiota.

15

16 Philipp Engel and Sudip Das jointly supervised the work.

17 #Address correspondence to Sudip Das, [sudip.das@unibe.ch](mailto:sudip.das@unibe.ch)

18 \* Present address: Slipa Kanungo, Instituto Gulbenkian de Ciência, Oeiras, Portugal.

19 Germán-Bonilla Rosso, Bioinformatics and Proteogenomics, Agroscope, Zürich,  
20 Switzerland

21

22

23

24 **Abstract**

25 Streptococci are one of the predominant and the most diverse genus in the human  
26 lung. Previously, we isolated human distal lung Streptococci from bronchoalveolar  
27 lavage fluid (BALF) as part of the human Lung Microbiota culture Collection (LuMiCol).  
28 Here, we performed whole genome sequencing, comparative phylogenomics and  
29 phenotypic characterization of six Streptococcal isolates representing the  
30 phylogenetic diversity of the genus in distal human lung. Here, we report five new  
31 species and one new subspecies including phylogenetic intermediates of commonly  
32 found Streptococci not limited to human lung. Pangenome analysis reveals gene  
33 content, evolutionary relationships, and metabolic functions shedding light on  
34 contribution of these Streptococci to lung microbial metabolism. Antimicrobial  
35 resistance gene analysis followed by MIC determination revealed macrolide,  
36 lincosamide and tetracycline resistance in lung Streptococci. We show the presence  
37 of capsular genes in lung streptococci both matching to the prototypical capsular  
38 genes (*cps*) and unique genes. Interestingly, the new *Streptococcus* isolate sp. nov.  
39 P2E5, genetically identical to the most prevalent *Streptococcus* in the human distal  
40 lung was revealed to be a phylogenetic intermediate between the *S. mitis* group and  
41 *S. pneumoniae*. It also harbors the pneumolysin (*ply*) gene and was found to have the  
42 serotype 21E. Finally, core genome phylogeny reveals that lung Streptococci are  
43 evolutionary distinct from oral Streptococcal isolates in expanded Human Oral  
44 Microbiome Database (eHOMD). Hence, these findings we reveal new  
45 phylogenetically distinct Streptococcal species from the human distal lung microbiota  
46 and its genetic diversity and metabolism to understand the microbial ecology of human  
47 lung.

48

## 49 **Importance**

50 A healthy human distal lung harbour characteristic microbial communities mostly  
51 composed of oropharyngeal taxa, which are facultative or obligative anaerobes despite  
52 lung being the medium of oxygen intake. However, little is known about the genetic  
53 and functional diversity of these bacteria owing to the lack of resources including  
54 availability of primary lung isolate from human samples. Therefore, we have  
55 established a large bacterial collection that covers all major phyla by cultivating human  
56 bronchoalveolar lavage fluid (BALF) under various conditions. *Streptococcus* is the  
57 most prevalent and diverse genera in the human lung microbiota. Using genetic and  
58 biochemical approaches, we studied six diverse lung isolates from our collection  
59 representing the actual Streptococcal diversity and identify these as new species and  
60 subspecies. We hypothesize that learning about the phylogenetic genetic diversity,  
61 preferred metabolism and molecular structures of these Streptococci will provide with  
62 new insights on the understudied microbial ecosystem of the human lung.

63

64

65 **Keywords:**

66 Human lung microbiota, viridans group streptococci (VGS), pangenome, human oral  
67 microbiota, novel streptococci

68

69

70

71

## 72 **Introduction**

73 The development of culture independent high-throughput DNA sequencing (both  
74 marker-gene amplicon sequencing and shotgun metagenomics) has made it possible  
75 to study the composition, diversity, and function of human microbial communities.  
76 Using culture-dependent and independent techniques we and others have shown that  
77 healthy human lungs harbour characteristic microbial communities(1–5). The lung  
78 microbiota is a complex and dynamic ecosystem composed of a diverse community  
79 of microorganisms, including bacteria, viruses, and fungi. The most prevalent bacterial  
80 phyla in lung are *Bacteroidetes* and *Firmicutes*, with low numbers of *Proteobacteria*  
81 and *Actinobacteria*(6). In addition, the biomass in the lung is relatively low compared  
82 to gut content with different genus level composition(7, 8). Most lung bacterial  
83 commensals have been shown to be of oral or supraglottic in origin(9, 10). However,  
84 the structure and composition of the lung bacterial communities are distinct(11). These  
85 differences may occur due to different oxygen conditions, pressure, pH, nutrients and  
86 distinct immune cell populations like airway macrophages that bacteria encounter  
87 during colonization(12). This suggests that the microbial ecology of lung and the  
88 interaction with the immune system is distinct from other sites on the human body(7,  
89 12). Despite the implications of lung-associated bacteria in lung health(13–16), our  
90 current understanding of resident lung microbiota is poor. (2)

91 With the advent of next generation sequencing came the ease of having taxonomic  
92 snapshot of a particular microbial niche leading to bacterial culture being overlooked.  
93 However, this is changing now with a broader realization of the importance of microbial  
94 cultivation and genotyping (17). In line with this, we performed large-scale culturing  
95 efforts with human BALF samples to obtain more than 300 bacterial isolates from 47  
96 species to build the open source bacterial biobank called [LuMiCol \(Lung Microbiota](#)

97 [culture Collection](#), Figure S1A). This covers the most prevalent species in human lung  
98 as well as important pathogens, observed by amplicon sequencing. This is an  
99 important resource that will facilitate experimental work on the human lung microbiota.  
100 We have also demonstrated that *Streptococci* were the most prevalent and diverse  
101 genus within the balanced pneumotype supporting homeostasis(18). Streptococci  
102 were amongst the most prevalent OTUs (5 out of 22, 16S rRNA gene identity) in  
103 bronchoalveolar lavage fluids, which was also apparent when cultivated to establish  
104 our bacterial collection. Our culture collection harbored six representative isolates, five  
105 matched to the most prevalent and abundant Streptococci (OTU\_11: P2E5, OTU\_34:  
106 P2D11, OTU\_42: P3D4, OTU\_57: P3B4, 369.3: OTU\_69) in human lung within our  
107 cohort (>97% 16S rRNA gene identity). Additionally, one isolate represented a rare  
108 *Streptococcus*.

109 Due to this diversity, we hypothesize that these Streptococci may represent the major  
110 metabolic pathways and provide valuable insights into the microbial ecology of the  
111 human distal lung. In this study, we sequenced the genomes of six *Streptococcus*  
112 isolates that represent each phylotype (97% 16S rRNA identity) followed by whole  
113 genome phylogenetic analysis, biochemical and metabolic tests. By doing these, we  
114 identified five new species and one novel subspecies of *Streptococcus*. Next, by  
115 employing pangenome analysis, we reveal orthologous gene content. Using a custom  
116 pipeline, we predicted metabolic pathways and macromolecular structures. We also  
117 reveal antibiotic resistance and virulence factors in these commensal streptococci.  
118 Finally, we compared lung streptococcal isolates with closely related genomes from  
119 expanded human oral microbiome database (eHOMD) to reveal their phylogenetic  
120 relationships.

121

## 122 **Results**

### 123 **Whole genome phylogenomics and phenotypic characterisation identifies new** 124 **streptococcal species and subspecies from distal human lung.**

125 We performed whole genome sequencing to obtain draft genomes of six lung isolates  
126 cultivated from BALF that represented the streptococcal diversity, using short read  
127 sequencing on the Illumina platform (Dataset S1, S2, S3). For species identification,  
128 we used both phylogenetic and biochemical approaches. Firstly, we performed  
129 genome-based identification using digital DNA-DNA Hybridization (dDDH)(19) on  
130 Type Strain Genome Server (TYGS) webtool(20). From this, we obtained species-  
131 level identification, the reference genomes of closest type strains and an outgroup  
132 taxon spanning a wide phylogenetic diversity (Table S1, S2, Dataset S4, S5).  
133 Additionally, we also obtained Genome BLAST Distance Phylogeny (GBDP)-based  
134 and full-length 16S rRNA-based phylogeny (Figure 1A, Figure S1B). Secondly, we  
135 performed pairwise whole genome Average Nucleotide Identity (FastANI(21))  
136 including the reference genomes and outgroup taxon (Figure 1B, Figure S2A). Thirdly,  
137 for further precision, we generated single-copy core gene phylogeny (Figure 2A,  
138 Figure S2B). Fourthly, we identified lung isolates using routine clinical microbiology,  
139 which included optochin resistance test, Matrix-Assisted Laser Ionization Time-Of-  
140 Flight (MALDI-TOF)-based protein spectral analysis and hemolysis (Table S3, Figure  
141 S2B, C). Finally, we performed standardized biochemical and metabolic panel tests  
142 (Strep API20, Biomerieux) for phenotypic characterization (Table S4, Figure S3). All  
143 lung isolates were confirmed to be Viridans Group Streptococci (VGS, Figure S2B,  
144 C)(22). Phylogenetic analysis suggested that five out of six isolates represent potential  
145 new species. (i) *Streptococcus* isolate sp. nov. P2E5 is a novel species within the *S.*  
146 *mitis* group, occupying an intermediate phylogenetic position between the *S. mitis* (ANI

147 93.26%) and *S. pneumoniae* (ANI 92.72%) clades (Figure 1A, B, 2A). This isolate  
148 exhibits a typical  $\alpha$ -hemolysis, proteomic analysis identified it as *S. mitis* whereas  
149 biochemical analysis indicated to *Gemella haemolysans*. (ii) *Streptococcus* isolate sp.  
150 nov. P2D11 is a new species within the *S. salivarius* group(23, 24). Interestingly, this  
151 isolate does not display any hemolysis ( $\gamma$ -hemolysis) and generated a unique  
152 biochemical pattern unlike a typical *Streptococcus* (Table S3, S4). (iii) *Streptococcus*  
153 isolate sp. nov. 369.3 is novel species that phylogenetically similar to *S. bovis* group  
154 and exhibits weak  $\alpha$ -hemolysis (Figure 1, 2, S2, Table S3, S4). However, proteomic  
155 analysis indicated *S. parasanguinis* and *S. australis* (Table S3) and biochemical  
156 identification show similarities to *Gemella morbillorum* with low discriminatory power  
157 (Table S4). (iv, v) *Streptococcus* isolates sp. nov. P3B4 and P3D4 represent two  
158 closely related new species with phylogenetic and spectral similarities to *S.*  
159 *parasanguinis*, displaying typical  $\alpha$ -hemolysis (P3B4 exhibited weak hemolysis).  
160 However, these two isolates exhibit biochemical and metabolic characteristics similar  
161 to *S. mitis* group. Lastly, *Streptococcus* P3E5 was identified as a new subspecies of  
162 *S. constellatus*, characterized by  $\beta$ -hemolysis, indicative of this specific this species.  
163 We henceforth, referred it as *S. constellatus* spp. nov. P3E5.

164

### 165 **Pangenome analysis sheds light on gene content and evolutionary relationship** 166 **amongst lung streptococci**

167 To understand the gene content of human lung *Streptococci* when compared to the  
168 reference genomes, we combined total proteins from 47 reference bacteria (Dataset  
169 S6) and 6 isolates (Dataset S7) to perform orthologous gene analysis using  
170 OrthoFinder(25). Overall in 53 genomes, we found 5232 orthogroups accounted for  
171 97.7% (99,280/101,643) of all proteins along with 150 strain-specific orthogroups

172 (2.8%) (Figure 2A, Dataset S8). Core genes represented 493 orthogroups (9.42%)  
173 with 315 single-copy core genes, which we used to construct a maximum likelihood  
174 evolutionary tree (Figure 2A) that corroborated the whole genome-based phylogeny  
175 (Figure 1). Next, we investigated the gene content of the Streptococcal isolates in  
176 comparison to each other. We combined all proteins from 47 reference genomes to  
177 construct a custom *Streptococcus* pan-proteome (Pan-Strep) database (Dataset S9)  
178 and compared each isolate to this using OrthoVenn2(26)(Dataset S10). The Pan-Strep  
179 database contained 31726 orthologous (169,727 proteins) with majority of the proteins  
180 present in lung isolates with few exceptions. *Streptococcus* isolate sp. nov. P2E5,  
181 P2D11, P3B4, P3D4, 369.3 and *S. constellatus* spp. nov. P3E5 contained 1970, 2086,  
182 1988, 2024, 1833 and 1891 genes belonging to 1780, 1846, 1872, 1883, 1757 and  
183 1696 gene clusters respectively. Furthermore, these isolates shared 977 genes  
184 including 957 single-copy core genes (Figure 2B). Interestingly, all isolates contained  
185 only 1 unique gene cluster each with none found in *Streptococcus* isolate sp. nov.  
186 P2E5. For functional categorisation, we performed Clusters of Orthologous Groups  
187 (COG) analysis in individual isolates and shared genes using eggNOG mapper(27).  
188 We found 22 COGs including 20 with known functional groups and 2 with unknown  
189 function (Figure 2C, Table S5, Dataset S11). Majority of the genes belonged to the  
190 COG category of unknown function (S) followed by translation, ribosomal structure  
191 and biogenesis (J), transcription (K) and Replication, recombination and repair (L).  
192 The most abundant cellular process was cell envelope biogenesis (M), and the most  
193 abundant metabolic genes were responsible for amino acid metabolism (E),  
194 carbohydrate metabolism (G) and inorganic ion metabolism (P). Contrastingly, cell  
195 motility (N), RNA processing and modification (A), extracellular structure (W),  
196 intracellular trafficking, secretion and vesicular transport (U) and lipid transport and



197 metabolism (I) were not prevalent. Although phylogenetically different, shared most of  
198 the COGs indicating similarities in basic cellular and metabolic functions.

199

200 **Metabolic and functional analysis of distal lung streptococci provide insights**  
201 **on the lung microbial ecosystem.**

202 Next, we predicted metabolic functions and macromolecular machineries using a  
203 custom rule-based based pipeline that included the GapMind, dbCAN and  
204 MacSysFinder tools (28–32) to comprehensively investigate the common catabolic  
205 and biosynthetic pathways, secretion systems, bacterial competence and  
206 carbohydrate-active enzymes (CAZymes)(33)(Figure 3, Table S6). The most  
207 prevalent mechanism for carbon catabolism in all streptococci including the lung  
208 isolates was the Embden-Meyerhof-Parnas (EMP) pathway. This was followed by  
209 Pentose Phosphate pathway (PPP) and Entner-Doudoroff (ED) pathway amongst  
210 majority of streptococci we tested. However, in majority of the bacteria, we observed  
211 either complete absence or incomplete canonical TCA cycle and oxidative  
212 phosphorylation (OXPHOS).

213 All lung isolates except *Streptococcus* isolate sp. nov. P2D11 possessed complete  
214 PPP, which was consistent within the *S. salivarius* group. We further investigated the  
215 number of carbohydrate-active enzymes (CAZymes, Figure 3, Figure S4A, Table S6,  
216 Dataset S12) and individual capacity to ferment sugars in culture (Figure S3). In total,  
217 the Pan-Strep database contained 5705 CAZymes subdivided into 6 families (Figure  
218 S4A). Compared to *S. pneumoniae* (108 CAZymes) that could ferment D-lactose, D-  
219 raffinose and D-trehalose. *Streptococcus* isolate sp. nov. P2E5 (mitis group, 74  
220 CAZymes) and *Streptococcus* isolate sp. nov. 369.3 (*S. bovis* group, 85 CAZymes)  
221 exhibited no sugar fermentation despite having GH1, 4 family of Glycosyl Hydrolases

222 (Figure S4B). *Streptococcus* isolate sp. nov. P3B4 (88 CAZymes) and P3D4 (87  
223 CAZymes) have a similar biochemical and metabolic profiles with D-lactose and D-  
224 raffinose fermentation capability, with the latter additionally fermenting D-sorbitol.  
225 *Streptococcus* isolate sp. nov. P2D11 (83 CAZymes) ferments D-mannose, D-sorbitol  
226 and D-lactose. *S. constellatus* novel. spp. P3E5 (61 CAZymes) was only able to  
227 ferment D-trehalose.

228 Next, we show that all isolates possessed pathways for the biosynthesis of most amino  
229 acids with a few exceptions (Figure 3, Table S6). All bacteria possess chorismate  
230 biosynthesis pathway (*aroG*, *aroB*, *aroD*, *aroE*, *aroL*, *aroA*, *aroC*), which can serve as  
231 an intermediate for biosynthesis of essential amino acids. In addition, all streptococci  
232 had complete pathways for nucleotide biosynthesis (purines and pyrimidines).  
233 However, most streptococci including lung isolates lacked the genes necessary for  
234 biosynthesis of electron acceptors and mediators such as Heme, NAD<sup>+</sup>, Coenzyme A  
235 and vitamin biosynthesis with interesting exceptions. Unlike *Streptococcus* isolate sp.  
236 nov. P3B4, *Streptococcus* isolate sp. nov. P3D4 can synthesize Vitamin B12. In  
237 addition, *Streptococcus* isolate sp. nov. P2D11 possessed the capability for Vitamin  
238 B6 biosynthesis consistent with the *S. salivarius* group.

239 We also investigated the presence of macromolecular structures such as secretion  
240 systems, diversity of competence and DNA uptake complexes in human lung  
241 commensal Streptococci. As expected, the most prevalent of these multi-protein  
242 complexes found in streptococci were the competence (*com*) proteins (Figure S5)  
243 responsible for natural competence, DNA uptake and transformation. These included  
244 competence stimulation peptides (CSP) and export protein (*ComB*, *ComC*)(34, 35)  
245 and the major response regulator *ComX*(36). Two types of DNA uptake complexes  
246 were found across all genomes: *ComE* proteins (*comEA*, *EB*, *EC*)(37) and the *ComF*

247 proteins (38) and the ComG pilus-like proteins (comGA, GB, GC, GD, GE, GF,  
248 GG)(39, 40).

249 Although no complete secretion systems were present in all genomes (Figure 3, Figure  
250 S5, Table S6), we still found Type IV secretion system (T4SS) proteins involved in  
251 conjugation(41). More specifically, proteins that we considered mandatory for a  
252 functional conjugative process were the ATPase complex system (VirB4)(42),  
253 coupling proteins (T4CP1, T4CP2)(41) and the type 4 toxin co-regulated pilus (TCP)  
254 subunit system (TcpA)(43) were present in all lung isolates. The accessory system  
255 including relaxases (MOBs)(44) were more variable across genomes with 5 MOB  
256 families (MOB<sub>B</sub>, MOB<sub>C</sub>, MOB<sub>Q</sub>, MOB<sub>T</sub>, MOB<sub>V</sub>) detected (Figure S5). Hence, we  
257 classified these as conjugation system proteins (CONJ)(28).

258

### 259 **Occurrence and prevalence of antimicrobial resistance and virulence factors in** 260 **human distal lung streptococci.**

261 Here, we investigated the presence of antimicrobial resistance (AMR) genes and  
262 virulence factors in the lung *Streptococcus* isolates using the ABRicate tool(45) (Figure  
263 4A, Table S7, S8). We also corroborated this with antibiotic susceptibility assays for  
264 all six isolates following EUCAST protocols that includes both disk diffusion assays  
265 and MIC tests (Table S2).

266 Comparison with the MEGARes database(46) revealed AMR genes in 27 reference  
267 genomes and 4 isolates (*Streptococcus* isolates isolate sp. nov. 369.3, P3D4, P3B4  
268 and P2E5), including multiple variants and /or copy numbers conferring resistance to  
269 6 classes of antimicrobials (Figure 4A). The most prevalent AMR was against  
270 Macrolides (66.7%; 18/27 genomes) followed by Tetracyclines (63%; 17/27) and  
271 Fluoroquinolones (37%; 10/27 genomes). Macrolide resistance genes (*mefA* and a

272 single copy of *msrD*(47, 48) were observed in 3 isolates apart from *Streptococcus*  
273 isolate sp. nov. P3D4. This was confirmed by antibiotic susceptibility assays (Table  
274 S2) where *Streptococcus* isolate sp. nov. P3B4 was resistant to Erythromycin (11 mm,  
275 MIC breakpoint = 4 mg/L), whereas *Streptococcus* isolate sp. nov. P3D4 was sensitive  
276 (27 mm). Although we didn't find any Lincosamide resistance genes (*IncC*, *IsaC*) in  
277 lung isolates, we tested susceptibility towards Clindamycin and performed D-test to  
278 distinguish between M- and MLS<sub>B</sub>-phenotype of macrolide resistance(48, 49). All  
279 isolates were sensitive to Clindamycin and showed no MLS<sub>B</sub>-phenotype, indicating  
280 only the presence of M-phenotype in human distal lung streptococci.

281 Tetracycline resistance genes were found in 4 isolates with *tetA*46 and *tetB*46(50)  
282 present in 3 isolates i.e., *Streptococcus* isolates sp. nov. 369.3, P3B4 and P3D4 and  
283 only one copy of *tetM* in one isolate i.e., *Streptococcus* isolate sp. nov. P2E5.  
284 Interestingly, all four isolates showed marginal resistance in diffusion assays and none  
285 in MIC tests. In addition, *Streptococcus* isolate sp. nov. P2D11 and P3E5 neither  
286 harboured the genes nor exhibited resistance phenotype. Interestingly, although we  
287 didn't find genes related to beta-lactam resistance, all lung isolates exhibited Oxacillin  
288 resistance (disk diffusion) and 3/5 isolates were resistant to Benzylpenicillin (MIC  
289 breakpoint = 0.25 mg/L, Table S2).

290 Virulence factor analysis using the VFDB(51) database revealed 20 genomes (2  
291 isolates and 18 references) harboring diverse virulence-related genes (Figure 4A). The  
292 highest number were observed in *S. pneumoniae* with *psaA* encoding for  
293 pneumococcal surface adhesin A, which plays a role in general and localized infection  
294 with *Streptococcus*(52, 53) the most prevalent across genomes. In one of the most  
295 interesting findings , we found pneumolysin (*ply*) gene, autolysin-encoding gene (*lytA*)  
296 and pneumococcal surface adhesin A (*psaA*) in the novel isolate *Streptococcus* P2E5,

297 the closest match to the most prevalent *Streptococcus* in human lung(18) and *psaA* in  
298 *Streptococcus* isolate sp. nov. P3E5. These genes are generally used to identify *S.*  
299 *pneumoniae*(54). Upon phylogenetic analysis, we show that *Streptococcus* isolate sp.  
300 nov. P2E5 Ply protein is similar to that of *S. pneumoniae* CCUG 28588 (98.3% identity)  
301 and *S. pseudopneumoniae* CCUG 49455 (98.7% identity) and type strain *S.*  
302 *pneumoniae* D39V (98.5% identity) (Figure 4B, Dataset S13). These results indicate  
303 the general abundance of AMRs but scarcity of virulence genes in VGS, including the  
304 lung isolates. Therefore, *Streptococcus* isolate sp. Nov. P2E5 is not only  
305 phylogenetically intermediate to *S. pneumoniae* and *S. mitis* but also in terms of  
306 virulence factors.

307

### 308 **Lung streptococci exhibit variable capsular diversity.**

309 Capsular polysaccharides common in commensal streptococci(55). However, we  
310 didn't observe all capsule genes in our virulence factor analysis. Hence, we  
311 investigated the presence of capsule genes in all genomes by protein BLAST against  
312 our Pan-Strep database using the prototypical *S. pneumoniae* D39 (Sp D39, serotype  
313 2) capsule (*cps*) locus i.e., the 17 genes located in between the *dexB* and *aliA* genes,  
314 as the query sequence followed by phylogenetic analysis of the matching genes  
315 (Figure 5A, Dataset S14). In addition, we also performed serotyping of the lung  
316 isolates by Quellung's test(56). Out of 53 genomes, 36 contained one or more capsular  
317 genes (Figure 5B, C) with the lowest (2 proteins) found in *Streptococcus*  
318 *pseudopneumoniae* ATCC BAA-960 and highest (22 proteins) in *Streptococcus*  
319 *salivarius* NCTC 8618. All six isolates possessed capsular proteins, which were a  
320 subset of Sp D39 (17 proteins) *cps* genes; P2D11 (12 proteins), P2E5 (13 proteins),  
321 P3E5 (11 proteins), P3B4 (12 proteins), P3D4 (11 proteins) and 369.3 (13 proteins).

322 The Sp D39 genes absent in isolates encoded for GTB-type glycotransferase  
323 superfamily of proteins(57) (Cps2G and CpsI), Capsular synthesis protein(58) Cps2H  
324 (*cps2H*), MATE-family protein(59) Cps2J (*cps2J*) and UDP-glucose 6-dehydrogenase  
325 Cps2K(60) (*cps2K*). However, unique capsular genes were also found in  
326 *Streptococcus* isolate sp. nov. P2D11 (hypothetical protein glycotransferase 1 family,  
327 protein ID EKHPBGBN\_01095), *Streptococcus* isolate sp. nov. 369.3  
328 (diaminopimelate decarboxylase; *lysA*, UDP-galactopyranose mutase; *glf*) and  
329 *Streptococcus* isolate sp. nov. P2E5 (UDP-galactopyranose mutase; *glf2*, UTP-  
330 glucose-1 phosphate uridylyltransferase; *cugP*). Interestingly, only *Streptococcus*  
331 isolate sp. nov. P2E5 tested positive for Quellung's test and was characterized to be  
332 serotype 21E. Hence, we compared its capsule genes with that of *S. pneumoniae*  
333 546/62 (Sp 546/62, reference for serotype 21) along with Sp D39 (serotype 2). We  
334 show that 11/13 capsule genes in P2E5 were high similarity to both Sp 546/62 and Sp  
335 D39 (Figure S6), indicating similarity to both serotypes.

336

### 337 **Core genome phylogeny reveals evolutionary relationship between lung, oral** 338 **and type strains of Streptococci.**

339 Although being of oral and supraglottic origin, the distal lung microbiota distinct (9–  
340 11). However, there is no study showing comparison at whole genome level. This  
341 prompted us to investigate the genetic proximity of the Streptococcal isolates  
342 cultivated from BALF and reference genomes from TYGS to that of the human oral  
343 microbiome. We used full-length 16S rRNA genes from the isolates and performed  
344 BLASTN against all genomes in the expanded Human Oral Microbiome Database  
345 (eHOMD)(61). This resulted in 47 representative genomes (best hits; > 97%16S rRNA  
346 gene identity), which we combined with lung streptococci and TYGS reference

347 genomes to perform core genome phylogeny (Figure 6, Table S9, Dataset S15).  
348 Overall, we did not observe body-site dependent pattern emerging rather all  
349 Streptococci were phylogenetically distributed regardless of the origin. Remarkably,  
350 lung streptococci stood out as phylogenetically distinct with one exception (Figure 6).  
351 *Streptococcus* isolate sp. nov. P2E5 was phylogenetically distinct with no closely  
352 related bacteria in the oral repertoire (Figure S7A). *Streptococcus* isolate sp. nov.  
353 P2D11 was closely related and intermediate to oral *S. salivarius* and *S. vestibularis*  
354 genomes in HOMD (Figure S7B). Interestingly, *S. constellatus* spp. nov. P3E5 was  
355 phylogenetically closer to oral *S. intermedius* but still within the *Streptococcus*  
356 *anginosus* group (Figure S7C). The phylogenetic placement of *Streptococcus* isolate  
357 sp. nov. 369.3 didn't change in relation to oral streptococci (Figure S7C). Finally,  
358 *Streptococcus* isolate sp. nov. P3B4 and P3D4 were observed to be phylogenetically  
359 distinct from both reference genomes and oral streptococci (Figure S7D). These  
360 results strengthen our findings of novel Streptococci and supporting the claims that  
361 lung microbiota is phylogenetically distinct from oral microbiota.

362

## 363 **Discussion**

364 Microbiota of the healthy lung is primarily derived from the oral and supraglottic  
365 niche(5, 8, 9, 12, 62, 63). This is also reflected in the lung post-transplant, where the  
366 oral taxa-dominant microbiota profile was associated with normal lung function and  
367 homeostasis(18, 64). Amongst all, Streptococci are the most spatiotemporally  
368 ubiquitous in oropharyngeal niche, upper and lower respiratory tract in healthy lung  
369 and allografts (5, 9, 11, 62–66).  
370 Likewise, in our previous study we have established an important resource called  
371 LuMiCol containing several lung bacterial isolates that match top lung taxa revealed

372 in amplicon sequencing(18). We also showed that *Streptococcus* is the most  
373 phylogenetically diverse and abundant genus in distal lung microbiota. However, due  
374 to limited resolution from amplicon sequencing, deeper genetic diversity in terms of  
375 specific species or strains were not known. Here, using robust phylogenomic analysis,  
376 comparative genomics and *Streptococcus*-specific phenotyping, we characterized 6  
377 different novel streptococcal isolates (Figure 1, 2, S3), which belonged to the highly  
378 heterogenous VGS and are evolutionary intermediates to already existing human-  
379 associated Streptococci, including both commensal and pathogenic species. We also  
380 categorized these into species groups whenever possible, which can be inconsistent  
381 (67). For example, *Streptococcus* isolate sp. nov. 369.3 is genetically related to *S.*  
382 *bovis* group II/1 (mannitol negative and beta-glucuronidase negative, Figure1, 2, Table  
383 S4) but shows biochemical similarity to the Nutritionally Variant Streptococcus (NVS)  
384 *G. morbillorum* (68). In addition, *Streptococcus* isolate sp. nov. P2E5 has both  
385 phenotypic features of the *S. mitis* group (Table S3) and biochemical features of *G.*  
386 *haemolysans* (Table S4). These observations along with the fact these novel isolates  
387 possess known orthologs upon comparison to the Pan-Strep database indicates intra-  
388 genera rather than an inter-genera gene transfer. Despite being high-quality these are  
389 not closed genomes and information on complex genetic structures might be missing  
390 especially considering the high genetic variation in *Streptococcus*. Hence, a  
391 combinatorial approach using short- and long-read sequencing should be the next  
392 appropriate step.

393 Human lung microbiota is primarily composed of facultative or obligate anaerobes,  
394 including the streptococci reported in this study (18). However, little is known about  
395 the microbial metabolism in the deep lung. Streptococci not only represent a larger  
396 subset of resident bacteria but are also temporally and spatially the most prevalent



397 genus(11, 18). Hence, its metabolic capabilities can provide crucial information on  
398 common catabolic and biosynthetic pathways within the lung microbiota. Although  
399 several genes for utilization and transport of sugars were present, there was a lack of  
400 canonical TCA cycle (Figure 3, Table S5). Additionally, the presence acetyl-CoA -  
401 Pyruvate/Lactate interconversion pathway (*ackA*, *pta* and *ldh*) (69) indicate a  
402 preference for anaerobic metabolism, which is in line with low glucose availability in  
403 airway epithelia(70). Hence, there might be two plausible pathways: the acetate-driven  
404 alternative TCA cycle (71) or pyruvate fermentation(69, 72).

405 The presence of complete pathways for acetate metabolism indicates its central role  
406 in the lung environment, which may be contributed mainly by commensal Streptococci  
407 (69, 73). Additionally, it is also an important short-chain fatty acid with  
408 immunomodulatory function in host gut and lung(74, 75) and shown to enhance killing  
409 of major lung pathogen *S. pneumoniae* by macrophages(76). However, these were  
410 mostly predictions, and we still lack information on nutritional preferences, which  
411 should be shown large-scale growth analysis on individual carbon sources.

412 Macromolecular structures in bacteria perform important functions in interacting with  
413 its environment. We revealed the presence multi-protein complex systems involved in  
414 bacterial competence, extracellular DNA uptake. All lung isolates harbor complete  
415 pathways for natural competence i.e., Com proteins including the pheromone peptides  
416 and regulators responsible for natural competence and extracellular DNA uptake  
417 complexes: ComE, F and G proteins (Figure S5)(34–40). However, this should be  
418 supported by further experimental induction of competence followed by DNA  
419 uptake(77). Lung streptococci also possess conjugative abilities shown by the  
420 presence of the VirB4)(42), TcpA)(43), T4CP1 and T4CP2(41) and may exchange  
421 genetic material with other genera in the community acquiring new traits.

422 Antibiotic resistance and virulence factors were mostly found in the mitis group VGS  
423 and pneumococci (Figure 4). Previously studies showed that antibiotic resistance is  
424 widespread in VGS and other human associated streptococci(78). Amongst the lung  
425 isolates, *Streptococcus* isolate sp. nov. P2E5 had most number with 7 genes (Table  
426 S6) and resistance pattern similar to other members of *S. mitis* group. As previously  
427 described for tetracycline resistance in oral streptococci(50), *tetM* encoding for  
428 ribosome protecting proteins was more common in our lung isolates than *tetAB*  
429 encoding for efflux pumps. Remarkably, presence of these genes did not manifest into  
430 phenotype (Table S6) apart from *Streptococcus* isolate sp. nov. 369.3. This could be  
431 due the requirement of additional genes or a result of altered gene regulation. Hence,  
432 it is challenging to conclude due the limitation that Tetracycline resistance in VGS is  
433 ill-defined by EUCAST due to insufficient evidence. (75). Interestingly, majority of lung  
434 isolates were resistant to the narrow spectrum beta-lactam Oxacillin, although we  
435 couldn't report resistance genes. However, remains inconclusive without investigating  
436 penicillin binding proteins (PBPs), which are crucial in conferring beta-lactam  
437 resistance to streptococci(79). Also, pneumococcus-specific virulence factors such as  
438 *ply*, *lytA* and *psaA* was found in *Streptococcus* isolate sp. nov. P2E5, which tested  
439 positive of pneumococcal polysaccharide capsule (serotype 21E) (Figure 5). However,  
440 it had less genes when compared to known serotype 21 reference genome Sp 546/62  
441 (Figure S6). The unique proteins in P2E5 such as UDP-galactopyranose mutase (*glf2*)  
442 and UTP-glucose-1 phosphate uridylyltransferase (*cugP*) may contribute to the  
443 specific capsule biosynthesis. However, this requires further investigation using other  
444 tests like immunodiffusion test(80) and heterologous expression of these genes to  
445 confirm its contribution. Remarkably, this isolate matched to the most prevalent  
446 *Streptococcus* in the human lung, which is associated with good lung function and

447 immunological balance. Presence of pneumococcal capsule serotypes has been  
448 reported in VGS and other human associated streptococci(60, 81). But P2E5 is unique  
449 as it an evolutionary intermediate with clear features of pneumococcus and *S. mitis*  
450 (Figure 1A, B, 2A). This phenotypic diversity and intermediary features amongst lung  
451 streptococci along with the presence of functional machineries for horizontal gene  
452 transfer may indicate a high degree of genetic exchange in the lung microbial  
453 environment. Previous studies have shown despite finding its origin in oral niche, the  
454 structure and composition of lung microbiota is distinct. Here, we provide genome-  
455 level evidence for the first time and show that lung microbiota remains phylogenetically  
456 distinct when compared to oral isolates from eHOMD.

457 Hence, to our knowledge, this is the first study to genome sequence novel lung  
458 bacterial isolates and perform comparative genomics to reveal crucial genetic,  
459 metabolic and evolutionary information filling the knowledge gaps in the field of  
460 microbial ecology of the human distal lung.

461

## 462 **Materials and Methods**

### 463 **Sample collection and ethics**

464 Sampling was performed and anonymized as previously described(18). The sampling  
465 via bronchoscopy was performed on individuals post lung transplant. Bronchoalveolar  
466 lavage fluid (BALF) was cultivated at random on different media and at different  
467 oxygen conditions. This sampling was approved by the local ethics committee (“Com-  
468 mission cantonale (VD’ d’éthique de la recherche su’ l’être humain – CER-VD”,  
469 protocol number 2018-01818) with written informed consent.

470

### 471 **Data and code availability**

472 All sequencing raw data were submitted to NCBI Short Read Archive under the  
473 BioProject [PRJNA1001255](#). Individual isolates were submitted under different  
474 BioSamples i.e., SAMN36797456, SAMN36797455, SAMN36797454,  
475 SAMN36797453, SAMN36797452, SAMN36797451. Processed data and  
476 supplementary datasets were uploaded on zenodo under DOI  
477 [10.5281/zenodo.10220079](#). Processed data including metaQUAST files, FASTA  
478 sequences and annotation files. All codes and pipelines are available on the GitHub  
479 [https://github.com/slipa17/Whole-genome-sequencing-and-comparative-genomics-](https://github.com/slipa17/Whole-genome-sequencing-and-comparative-genomics-of-human-lung-streptococcal-isolates)  
480 [of-human-lung-streptococcal-isolates](#). Details on any scripts (.sh files), workflows (.md  
481 /.Rmd or .R files) and parameters (.txt files), which are mentioned throughout can be  
482 be found on the GitHub page.

483

#### 484 **Bacterial growth and media**

485 For routine cultivation, all Streptococci were grown for 24 – 48 hours on Columbia agar  
486 (Oxoid, UK) with 5% defibrinated sheep blood (Thermo Scientific, USA) at 37°C in  
487 presence of 5% CO<sub>2</sub> and 95 % relative humidity or in a vinyl anaerobic chamber with  
488 < 5 ppm O<sub>2</sub> (Coy labs, USA) at 35°C with moisture control. For broth cultures, one or  
489 two isolated colonies were picked and inoculated in a polypropylene culture tube (with  
490 cap) containing 2 ml of Todd-Hewitt Broth (Oxoid, UK) supplemented with 0.5% Yeast  
491 Extract (Oxoid, UK). The cultures were incubated under the same conditions as  
492 mentioned above and strictly without agitation.

493

#### 494 **Bacterial DNA isolation and genome sequencing**

495 Some streptococcal isolates exhibited unusual physical properties upon growth on  
496 semi-solid media, which included dry, flaky colony texture, difficulty in resuspension in

497 buffer and recalcitrance. Hence, a custom bacterial DNA isolation protocol was used.  
498 This process involved sequential lysis of bacteria using both enzymatic action and  
499 mechanical shearing followed by extraction using QIAamp DNA Mini Kit (QIAGEN,  
500 Germany). Bacteria were grown as broth cultures and harvested followed by  
501 resuspension in 200 µl of Gram-positive lysis buffer (20 mM Tris-HCl, pH 8.0, 2 mM  
502 EDTA, 1.2% Triton X-100) containing 1 mg/ml lysozyme and 100 µg/ml RNase A. The  
503 mixture was incubated for 30 min at 37°C with gentle agitation. After this, the volume  
504 was brought up to 500 µl with Gram-positive lysis buffer and the mixture was  
505 transferred to screw cap tubes in tubes containing 200 mg of 0.1-mm acid-washed  
506 zirconia beads and homogenized using a FastPrep-25 5G instrument (2 rounds of 30 s  
507 with the power set to 6), as previously described(28). This was followed by  
508 centrifugation at maximum speed for 5 minutes at room temperature. The debris-free  
509 supernatant was carried over to the QIAamp kit protocol, which involves incubation  
510 with Proteinase K followed by column-based extraction steps. Genomic libraries for  
511 Illumina sequencing libraries were prepared in-house using the Vazyme TruePrep  
512 DNA library preparation kit following manufacturer's instructions. Multiplexing was  
513 performed using Nextera i7 adaptors. Sequencing was performed on the Illumina  
514 HiSeq 2500 instrument at the Genomics Technology Facility, University of Lausanne,  
515 Switzerland using two simultaneous lanes for avoiding lane-bias generating 150 bp  
516 pair-end reads.

517

### 518 **Bacterial genome assembly and annotation**

519 Read quality control and trimming was performed using FastQC v0.11.9(82) and  
520 Trimmomatic v0.39(83) (parameters: PE -phred33 AIIllumina-Peadapters.fa:3:25:7  
521 LEADING:9 TRAILING:9 SLIDINGWINDOW:4:15 MINLEN:60). SPAdes(84) (-careful

522 option, v3.15.2) was used for *de novo* assembly of bacterial genomes using  
523 *run\_spades.sh* (Dataset S2). Different parameters were used assess the quality of  
524 assemblies (*spades\_param.txt*). Assemblies were evaluated for its quality and  
525 completeness using metaQUAST (Quality Assessment Tool for Genome Assemblies)  
526 v5.0.2(85) and checkM v1.0.13(86). All genome statistics and metaQUAST HTML  
527 report files were created for summaries of each assembly task (Dataset S1). Draft  
528 genome scaffolds were annotated using prokka v1.13(87) using *prokka.sh* (Dataset  
529 S3).

530

### 531 **Genome-based bacterial identification using Type Strain Genome Server (TYGS)**

532 For identification of closely related *Streptococcus* species and Genome-based  
533 Distance Phylogeny (GBDP), the DNA FASTA files of the isolates were submitted to  
534 [TYGS](#)(20) web portal. This tool uses for genome and 16S rRNA BLAST with clusters  
535 species and subspecies to identify species and report nearest neighbours. The output  
536 includes genome and 16S rRNA based phylogenetic trees, which can be exported.  
537 Visualization of these trees and associated metadata was done using [iTOL](#)(88) v6.8.1  
538 webtool.

539

### 540 **Data extraction from public databases**

541 All reference genomes from TYGS analysis (Dataset S4) including GenBank data files  
542 (.gbk, .gbff), annotation features (.gff), Nucleotide (.fna, .fa) and Protein (.faa) FASTA  
543 files were downloaded from NCBI FTP server using NCBI Datasets command-line  
544 tools (CLI), using *NCBI datasets and BLAST.md*. In case any genome assembly did  
545 not accompany translations or proper annotations, which were then annotated using  
546 PROKKA. Human oral isolates genomes were downloaded from [eHOMD Genomes](#).

547 For consistency and reproducibility of further analyses, all genomes were reannotated  
548 using prokka v1.13(87) using *prokka.sh* (Dataset S5, S14).

549

### 550 **Pairwise average nucleotide identity (ANI)**

551 Pairwise Average Nucleotide Identity (ANI) including visualization between  
552 *Streptococcus* genomes was performed using *FastANI.md* that combined ANIcluster  
553 and FastANI command line tools (89, 90). The output files contained query genome,  
554 reference genome, ANI values, count of bidirectional fragment mappings, and total  
555 query fragments.

556

### 557 **Pangenome analysis and orthologous group:**

558 This involved three steps: 1. Creation of a Pan-Strep database: all proteins FASTA  
559 files available from reference genomes provided by the TYGS analysis(20) were  
560 concatenated to form the pan-proteome database, which was used both for orthology  
561 and BLAST database (Dataset S9). 2. Comparison of isolates to Pan-Strep database:  
562 [OrthoVenn2](#) webtool(26) was used classify orthologous gene clusters in each isolate  
563 using Pan-Strep database as reference. The E-value and inflation value were set at 1  
564  $\times 10^{-5}$  and 1.5 respectively. The distribution of the shared (only between isolates) and  
565 unique orthologous clusters, total protein count was represented as a Venn diagrams  
566 (Dataset S10). 3. Comparison of all genomes to each other: All protein FASTA files  
567 from genomes were analyzed together using Orthofinder v2.5.5(25) to obtain the  
568 single copy core orthogroups. This provides multiple statistics on the genetic contents  
569 including orthologs, xenologs and shared genes and single-copy shared genes  
570 (Dataset S8).

571

## 572 **Single copy core genome phylogeny**

573 A list of single copy core (shared) genes (Dataset S7) used to extract corresponding  
574 protein FASTA from each genome and create individual FASTA files containing core  
575 genes using *Append\_concatenate\_extract\_grep\_Linux\_log.md* and  
576 *Extraction\_of\_headers\_fasta\_sequences\_and splitted\_files.md*. All proteins in each  
577 FASTA file were self-concatenated to create a single sequence FASTA with one  
578 header. All resulting sequences were aligned using MAFFT v7.52 command line  
579 tool(91) and Maximum-likelihood (ML) tree was computed using the FastTree  
580 v2.1.11(92) or RAxML v 8.2.12 command line tool (93). Visualization of these trees  
581 was done using [iTOL](#)(88) v6.8.1 webtool.

582

## 583 **Multiple sequence alignment and phylogeny**

584 All alignments (Protein or DNA FASTA sequences) were performed using MAFFT  
585 v7.52(91) command line tool and phylogenetic trees were constructed using FastTree  
586 v2.1.11 or RAxML v8.2.12 command line tool (93), unless and otherwise specified.  
587 The workflow *Alignment and phylogeny.md* was used that includes MAFFT for  
588 MacOSX (version 7.505) with the argumen--'auto' for automatic detection of  
589 parameters. The output FASTA file containing the aligned sequences was masked  
590 (50% sites stripped) using Geneious Prime software (New Zealand) and exported as  
591 PHYLIP format. Maximum-likelihood (ML) trees was computed with the alignment file  
592 using either RAxML v8.2.12 (model 'PROTGAMMAAUTO', 100 rapid bootstrapping)  
593 or FastTree v2.1.11 (LG + CAT substitution model for protein or default -nt model for  
594 nucleotides). Visualization of these trees was done using [iTOL](#)(88) v6.8.1 webtool.

595



## 596 **Comparison with genomes from expanded Human Oral Microbiome Database** 597 **(eHOMD)**

598 For this analysis, full-length 16S rRNA gene sequences of the isolates were used as  
599 queries to search on the [HOMD RefSeq BLAST Server](#) (perform BLASTN v2.12.0,  
600 HOMD\_16S\_rRNA\_RefSeq\_V15.23.p9 database). The hits (>97% identity and >99%  
601 coverage) were selected and its genomes were downloaded from eHOMD. These  
602 genomes were reannotated using prokka v1.13(87) using *prokka.sh*. Comparison with  
603 the lung isolates and reference type strains from TYGS was performed by single-copy  
604 core genome phylogeny (Dataset S14). The protein FASTA files all sources were used  
605 to run Orthofinder v2.5.5 to obtain the single copy core orthogroups. These orthogroup  
606 FASTA files were extracted using *Append\_concatenate\_extract\_grep\_Linux\_log.md*  
607 and *Extraction\_of\_headers\_fasta\_sequences\_and\_splitting\_files.md* that rearranges  
608 these into single-copy core proteins according to each genome. These protein files  
609 were then aligned with MAFFT v7.52 and tree was constructed using FastTree v2.1.11  
610 (LG + CAT substitution protein model). Tree visualization was done using [iTOL](#)(88)  
611 v6.8.1 webtool.

612

## 613 **Clinical microbiology, identification and phenotyping**

614 Species identification was done as described(94). In brief, routine bacterial  
615 identification was performed using a combination of MALDI-TOF, serotyping and  
616 functional assays at the Institute for Infectious Diseases (IFIK), University of Bern,  
617 Switzerland. Streptococci were grown on Columbia agar with 5% defibrinated sheep  
618 blood (Biomérieux) (CSBA plates) for 24 hours at 37°C in presence of 5% CO<sub>2</sub>. Single  
619 colonies were picked for MALDI-TOF analysis. In addition, bacterial were spread on  
620 CSBA plates as lawn cultures and Optochin disks (Sigma) was placed with a sterile

621 forceps. The plates were incubated for 24 hours at 37°C in presence of 5% CO<sub>2</sub> and  
622 inhibition zones were observed. A positive zone of inhibition (sensitive, 15 mm) is  
623 usually in case of pneumococci and no inhibition (resistant) is usually seen in viridans  
624 group or other alpha-hemolytic streptococci. The type of hemolysis was also noted  
625 during this assay. Finally, serotyping was performed using standard Quellung's  
626 (Neufeld) reaction(95) towards capsular polysaccharide as described(96).

627

### 628 **Antibiotic Susceptibility Tests**

629 Antibiograms and Minimum Inhibitory Concentration (MIC) studies were performed at  
630 the Institute of Infectious Diseases, University of Bern, Switzerland, according to the  
631 criteria established by European Committee on Antimicrobial Susceptibility Testing  
632 ([EUCAST](#))(96). Bacteria were grown on Muller-Hinton agar for Fastidious organisms  
633 (MH supplemented with 5% defibrinated horse blood and 20 mg/l β-NAD).  
634 Antibiograms were performed using disk diffusion method (disk content in µg),  
635 observing zone of inhibition (diameter in mm). MIC determination was performed using  
636 Etest® strips (Biomérieux, France, described in µg/ml). The values were tallied with  
637 the EUCAST v13.0 clinical breakpoint tables for interpretation of results. Macrolide-  
638 inducible resistance to clindamycin test (D-test) to test for assessing macrolide-  
639 lincosamide-streptogramin B (MLS<sub>B</sub>) resistance, was performed by placing  
640 Erythromycin and Clindamycin disks 12-20 mm apart (edge to edge). Appearance of  
641 antagonism (the D phenomenon) was observed to detect any inducible clindamycin  
642 resistance.

643

### 644 **Biochemical identification of *Streptococcaceae***

645 Streptococcal identification was carried out by using API® 20 Strep panel (Biomérieux,  
646 France), a standardized system combining 20 biochemical tests, according to  
647 manufacturer's instructions. Briefly, all isolates were first grown as previously  
648 described (Methods: Bacterial growth and media). Bacteria were collected using a  
649 sterile cotton swab and mixed in 2ml API® Suspension Medium to achieve turbidity  
650 more than McFarland standard 4 before distributing into the given cupules in the panel  
651 strips and incubating according to manufacturer's instructions. In this case, both 4-  
652 and 24-hour readings were performed as some isolates tend to exhibit delayed effects.  
653

## 654 **Functional annotation and analysis**

### 655 ***Metabolic profiling and macromolecular structures prediction***

656 Predictions of key metabolic pathways and macromolecular structures were  
657 performed by using a custom genome profiler as previously described(28). A set of  
658 rules were applied to conclude the presence and completeness of predicted pathways.  
659 These can be found in the scripts defining the rules at  
660 [https://github.com/g Barton/Publication\\_Sarton-Loheac\\_2022](https://github.com/g Barton/Publication_Sarton-Loheac_2022). Metabolic predictions  
661 were performed using GapMind(30) for amino acid biosynthesis and dbCAN(97) for  
662 Carbohydrate Active enzymes(98) ([CAZyme](#)) analysis (Dataset S11). For secretion  
663 systems and other macromolecular structures, MacSyFinder(99) was used. The  
664 results were visualized along as heatmap depicting system completeness (%) along  
665 with a single copy core genes-based phylogeny using [iTOL](#)(88) v6.8.1 webtool.  
666

### 667 ***COG categories***

668 eggNOG (evolutionary genealogy of genes: nonsupervised orthologous groups) was  
669 used to predict the functional annotation of genes. Representative genes from six

670 individual Streptococcal species and the 977 shared genes among the six isolates  
671 obtained from the OrthoVenn2 tool were used. Shared genes from any one of the files  
672 containing the protein FASTA sequences of the 6 isolates was performed using  
673 *Append\_concatenate\_extract\_grep\_Linux\_log.md* and  
674 *Extraction\_of\_headers\_fasta\_sequences\_and splitted\_files.md*. Protein (up to  
675 100,000) sequences from the genomes of individual species along with those shared  
676 977 extracted sequences was uploaded separately to the eggno-mapper v2 web  
677 tool(27). The individual output files (.xlsx) were downloaded (Dataset S10). The COGs  
678 assigned to different individual proteins was used to calculate frequencies in each of  
679 the isolates using *COG\_analysis\_Rscripts.md*. This was visualized as heatmap along  
680 with a single copy core genes-based phylogeny using [iTOL](#)(88) v6.8.1 webtool.

681

### 682 **Anti-microbial and Virulence genes analysis**

683 Genes coding for antimicrobial resistance (AMR) and known virulence factors were  
684 detected by *Antimicrobial\_resistance\_virulence\_gene\_search.md* that uses [ABRicate](#)  
685 [tool](#) (v 1.0.1). The databases used were MEGARes(46) and VFDB(51) for AMR and  
686 virulence factors respectively. Genomic DNA FASTA sequences of six Streptococcal  
687 species along with the 47 reference genomes were used as input files. The output text  
688 file (.csv) was summarized using `abrica -- summary` command, which included  
689 genome ID, gene name, percent coverage and number of genes found. The results  
690 were visualized along as heatmap depicting percentage (%) coverage along with a  
691 single copy core genes-based phylogeny using [iTOL](#)(88) v6.8.1 webtool.

692

### 693 **Screening of Streptococcal capsule in silico**

694 For capsule analysis, the prototypical *Streptococcus pneumoniae* capsule locus was  
695 used as the query and BLAST command line tool was used to search for matching  
696 proteins in the isolates and reference genomes. Protein FASTA sequences of capsule  
697 genes i.e., 17 proteins encoded by *cpsA-T* were extracted from the D39V genome and  
698 saved as a single FASTA file (Dataset S13). For creating a custom BLAST database,  
699 all proteins from 6 isolates and 47 reference genomes were combined to create the  
700 Pan-Strep database containing 80,364 protein sequences. The database was created  
701 under the 'ncbi-datasets' environment in command line using the 'makeblastdb'  
702 function. Full script and analysis steps are described in 'NCBI datasets and  
703 BLAST.md'. The FASTA file containing query capsule proteins was used to perform  
704 BLASTP 2.6.0 search and the output was obtained in a tabular format. A set of rules  
705 were decided for the protein to be considered as a significant match. These were: (i)  
706 E-value cut-off at 0.0, (ii) Bit score > 200 and (iii) identity threshold of > 50%. The  
707 protein IDs of the resulting matches used to extract the FASTA sequences using  
708 *Append\_concatenate\_extract\_grep\_Linux\_log.md* and  
709 *Extraction\_of\_headers\_fasta\_sequences\_and\_split\_files.md*. The capsule gene  
710 sequences were rearranged according to bacterial genomes and were concatenated  
711 followed by *Alignment and phylogeny.md* (50% sites stripped, RAxML model  
712 'PROTGAMMAAUTO', 100 rapid bootstrapping). The results were visualized along as  
713 barchart showing distribution and number of capsule genes found in each genome  
714 along phylogeny using [iTOL](#)(88) v6.8.1 webtool. Gene plots were constructed using  
715 pyGenomeViz(100) to visualize the genomic arrangement and synteny of lung  
716 Streptococcal isolates in comparison to the prototypical capsule locus of  
717 *Streptococcus pneumoniae* D39.

718

## 719 **Pneumolysin analysis**

720 Pneumolysin protein encoded by the *ply* gene, from the prototypical gene sequence  
721 from *Streptococcus pneumoniae* D39 was used as a query to perform BLASTP 2.6.0  
722 search against the custom Pan-Strep database. The hits found were then extracted  
723 from the respective genomes and followed by *Alignment and phylogeny.md*  
724 (MUSCLE(101), RAxML model 'PROTGAMMAAUTO', 100 rapid bootstrapping)  
725 (Dataset S12). For visualization, a heatmap with distance matrix was constructed in R  
726 and gene plot was made using pyGenomeViz.

727

## 728 **Gene plots and visualization of genomic features**

729 For visualizing sequence similarity and comparison of gene arrangements between  
730 multiple genomes pyGenomeViz-MMSeqs v0.3.2 tool was to plot genomic features.  
731 This tool was used in the conda environment with default parameters using Genbank  
732 files (.gb or .gbk) as input files. Output data contained the reciprocal best hits file (.tsv)  
733 and the syntenic plots between the genomes. In synteny plots, pairwise sequence  
734 similarity can be observed between genomes or coding sequences. The position of  
735 the genes can be compared using the genomic coordinates. This synteny analysis  
736 also informs about reciprocal mapping to identify regions of similarity and orthologous  
737 genes for understanding the evolutionary relationships.

738

## 739 **Acknowledgements**

740 Funding for this work came from several agencies for different contributors to this  
741 work, which are as follows. Sudip Das was funded by Marie Skłodowska-Curie  
742 Individual Fellowship ("HUMANITY", Grant no. 800301) and an Interdisciplinary grant  
743 from the Faculty of Biology and Medicine of the University of Lausanne (awarded to

744 Philipp Engel as supervisor, Grant no. 26075716). Philipp Engel is funded by an ERC  
745 StG (“MicroBeeOme”, Grant No. 714804), two Swiss National Science Foundation  
746 grants (Grant no. 31003A\_160345 and 31003A\_179487). Germán-Bonilla Rosso was  
747 funded by HFSP Young Investigator grant (awarded to Philipp Engel as supervisor,  
748 Grant no. RGY0077/2016). We thank Prof. Dr. Jan-Willem Veening and his team at  
749 the Department of Fundamental Microbiology, University of Lausanne, Lausanne,  
750 Switzerland for providing the bacterial strains *S. pneumoniae* D39 and *S. mitis* NCTC  
751 12261.

752

### 753 **Figure Legends**

754 **Figure 1. Whole genome phylogeny reveals novel streptococci from human**  
755 **distal lung microbiota.** Comparison of human distal lung streptococci (red gradient)  
756 to closely related type strains (gray gradient) in the TYGS database. **A.** Whole  
757 Genome BLAST Distance Phylogeny (GBDP) using FASTME, where colored boxes  
758 represent species and subspecies clusters, and blue-colored gradient boxes represent  
759 GC content (%). **B.** Heatmap showing pairwise Average Nucleotide Identity (ANI %)  
760 between human distal lung streptococci (red gradient) to closely related type strains  
761 (gray gradient). ANI values > 80% were colored as grey in the heatmap.

762

763 **Figure 2. Core genome phylogeny and orthologous gene content of human**  
764 **distal lung streptococci** **A.** Single copy core-genome phylogeny and orthologous  
765 gene analysis of lung streptococci (red gradient) and closely related type strains (gray  
766 gradient). Maximum-likelihood tree computed by FastTree on concatenated amino  
767 acid sequences of 315 single-copy core proteins using LG + CAT substitution model.  
768 Bar graphs show total number of genes (gray bars), percentage (%) genes in

769 orthogroups (green and red stacked bars), number genes that weren't assigned to  
770 orthogroups (blue bars), number of specific-specific orthogroups (dark brown) and  
771 number of genes in each (light brown). **B.** Venn diagram showing 977 shared proteins  
772 between lung streptococcal isolates generated using OrthoVenn2. **C.** Heatmap  
773 showing functional classification of proteins from lung streptococcal isolates using  
774 eggNOG depicted according to maximum-likelihood phylogeny computed by RAxML  
775 on concatenated amino acid sequences of 957 single-copy core proteins. Bar graphs  
776 show number of shared genes (y-axis) in each functional category (x-axis, colored  
777 strips).

778

779 **Figure 3. Function and metabolic prediction predictions of human distal lung**  
780 **streptococci.** Heatmap showing completeness of predicted pathways depicted  
781 according to maximum-likelihood phylogeny computed by FastTree on concatenated  
782 amino acid sequences of 315 single-copy core proteins using LG + CAT substitution  
783 model in individual human distal lung streptococci (red gradient) and closely related  
784 type strains (gray gradient). A custom rule-based pipeline was used to predict  
785 metabolic pathways; catabolic (red stripe), biosynthetic pathways (green stripe) and  
786 macromolecular systems (purple stripe). Stacked bar charts show counts (y-axis) of  
787 Carbohydrate-active enzymes (CAZymes) in individual genomes (x-axis).

788

789 **Figure 4. Prevalence of antimicrobial resistance and virulence factors in human**  
790 **distal lung streptococci.** **A.** Heatmaps depicting coverage of antimicrobial resistance  
791 (left panel) and virulence gene (right panel) revealed by ABRicate depicted according  
792 to maximum-likelihood phylogeny computed by FastTree on concatenated amino acid  
793 sequences of 315 single-copy core proteins using LG + CAT substitution model in



794 individual human distal lung streptococci (red gradient) and closely related type strains  
795 (gray gradient). Bar graphs show number of antimicrobial resistance (yellow) and  
796 virulence (purple) genes (y-axis) in each category (x-axis, colored strips). **B.** Heatmap  
797 depicting pairwise distance matrix after alignment generated using MUSCLE, showing  
798 percentage identities of Pneumolysin protein between *Streptococcus* isolate P2E5 and  
799 pneumococci. Maximum-likelihood phylogeny of Ply proteins computed by by RAxML  
800 and gene plot for genomic coordinates and synteny of *ply* genes generated by  
801 pyGenomeViz. This compared the percentage identity (shading depth) and relative  
802 gene orientation (grey for links with the same direction, pink for reverse direction). This  
803 compared the percentage identity (shading depth) and relative gene orientation (grey  
804 for links with the same direction, pink for reverse direction).

805

806 **Figure 5. Investigation of Streptococcal capsular polysaccharide synthesis and**  
807 **its relative genomic arrangement. A.** Workflow showing our capsular analysis  
808 pipeline from BLAST search to phylogeny. **B.** Maximum-likelihood phylogeny computed  
809 computed by by RAxML after alignment of capsular proteins found in individual human  
810 distal lung streptococci (red gradient) and closely related type strains (gray gradient).  
811 Bar graphs show number of capsule protein BLAST hits (y-axis) in individual genomes  
812 (x-axis). **C.** Comparison of capsular polysaccharide synthesis (*cps*) genes (colored  
813 arrows) in lung isolates to the prototypical *S. pneumoniae* D39 capsule operon. Gene  
814 plot generated by pyGenomeViz for comparing genomic coordinates and synteny with  
815 percentage identity (shading depth) and relative gene orientation (grey for links with  
816 the same direction, pink for reverse direction).

817

818 **Figure 6. Phylogenetic comparison of human oral and distal lung streptococci.**

819 Single copy core-genome phylogeny comparing lung streptococcal isolates cultivated  
820 from BAL (red gradient), Reference genomes from TYGS (gray gradient) and human  
821 oral streptococci from eHOMD (green gradient). Maximum-likelihood tree computed  
822 by FastTree on concatenated amino acid sequences of 26 single-copy core proteins  
823 using LG + CAT substitution model.

824

825

826 **References**

- 827 1. Charlson ES, Bittinger K, Haas AR, Fitzgerald AS, Frank I, Yadav A, Bushman FD,  
828 Collman RG. 2011. Topographical continuity of bacterial populations in the healthy  
829 human respiratory tract. *Am J Respir Crit Care Med*  
830 <https://doi.org/10.1164/rccm.201104-0655OC>.
- 831 2. Das S, Bernasconi E, Koutsokera A, Wurlod D-A, Tripathi V, Bonilla-Rosso G,  
832 Aubert J-D, Derkenne M-F, Mercier L, Pattaroni C, Rapin A, von Garnier C, Marsland  
833 BJ, Engel P, Nicod LP. 2021. A prevalent and culturable microbiota links ecological  
834 balance to clinical stability of the human lung after transplantation. *Nat Commun*  
835 12:2126-undefined.
- 836 3. Pattaroni C, Watzenboeck ML, Schneidegger S, Kieser S, Wong NC, Bernasconi E,  
837 Pernot J, Mercier L, Knapp S, Nicod LP, Marsland CP, Roth-Kleiner M, Marsland BJ.  
838 2018. Early-Life Formation of the Microbial and Immunological Environment of the  
839 Human Airways. *Cell Host Microbe* 24:857-865.e4.
- 840 4. Bernasconi E, Pattaroni C, Koutsokera A, Pison C, Kessler R, Benden C, Soccia PM,  
841 Magnan A, Aubert J-D, Marsland BJ, Nicod LP. 2016. Airway Microbiota Determines  
842 Innate Cell Inflammatory or Tissue Remodeling Profiles in Lung Transplantation. *Am*  
843 *J Respir Crit Care Med* 194:1252–1263.
- 844 5. Venkataraman A, Bassis CM, Beck JM, Young VB, Curtis JL, Huffnagle GB, Schmidt  
845 TM. 2015. Application of a Neutral Community Model To Assess Structuring of the  
846 Human Lung Microbiome. *mBio* 6.
- 847 6. Marsland BJ, Gollwitzer ES. 2014. Host–microorganism interactions in lung diseases.  
848 *Nat Rev Immunol* 14:827–835.
- 849 7. Man WH, De Steenhuijsen P, WAA, Bogaert D. 2017. The microbiota of the  
850 respiratory tract: Gatekeeper to respiratory health. *Nat Rev Microbiol*  
851 <https://doi.org/10.1038/nrmicro.2017.14>.
- 852 8. Dickson RP, Huffnagle GB. 2015. The Lung Microbiome: New Principles for  
853 Respiratory Bacteriology in Health and Disease. *PLoS Pathog* 11:e1004923.
- 854 9. Bassis CM, Erb-Downward JR, Dickson RP, Freeman CM, Schmidt TM, Young VB,  
855 Beck JM, Curtis JL, Huffnagle GB. 2015. Analysis of the Upper Respiratory Tract  
856 Microbiotas as the Source of the Lung and Gastric Microbiotas in Healthy Individuals.  
857 *mBio* 6.

- 858 10. Dickson RP, Erb-Downward JR, Martinez FJ, Huffnagle GB. 2016. The Microbiome  
859 and the Respiratory Tract. *Annu Rev Physiol* 78:481–504.
- 860 11. Dickson RP, Erb-Downward JR, Freeman CM, McCloskey L, Beck JM, Huffnagle  
861 GB, Curtis JL. 2015. Spatial Variation in the Healthy Human Lung Microbiome and  
862 the Adapted Island Model of Lung Biogeography. *Ann Am Thorac Soc* 12:821–830.
- 863 12. Lloyd CM, Marsland BJ. 2017. Lung Homeostasis: Influence of Age, Microbes, and  
864 the Immune System. *Immunity* 46:549–561.
- 865 13. Whelan FJ, Heirali AA, Rossi L, Rabin HR, Parkins MD, Surette MG. 2017.  
866 Longitudinal sampling of the lung microbiota in individuals with cystic fibrosis. *PLoS*  
867 *One* <https://doi.org/10.1371/journal.pone.0172811>.
- 868 14. Hilty M, Burke C, Pedro H, Cardenas P, Bush A, Bossley C, Davies J, Ervine A,  
869 Poulter L, Pachter L, Moffatt MF, Cookson WOC. 2010. Disordered microbial  
870 communities in asthmatic airways. *PLoS One*  
871 <https://doi.org/10.1371/journal.pone.0008578>.
- 872 15. Mika M, Nita I, Morf L, Qi W, Beyeler S, Bernasconi E, Marsland BJ, Ott SR, von  
873 Garnier C, Hilty M. 2018. Microbial and host immune factors as drivers of COPD.  
874 *ERJ Open Res* 4:00015–02018.
- 875 16. Woo TE, Lim R, Heirali AA, Acosta N, Rabin HR, Mody CH, Somayaji R, Surette  
876 MG, Sibley CD, Storey DG, Parkins MD. 2019. A longitudinal characterization of the  
877 Non-Cystic Fibrosis Bronchiectasis airway microbiome. *Sci Rep*  
878 <https://doi.org/10.1038/s41598-019-42862-y>.
- 879 17. Afrizal A, Hitch TCA, Viehof A, Treichel N, Riedel T, Abt B, Buhl EM, Kohlheyer D,  
880 Overmann J, Clavel T. 2022. Anaerobic single-cell dispensing facilitates the  
881 cultivation of human gut bacteria. *Environ Microbiol* [https://doi.org/10.1111/1462-](https://doi.org/10.1111/1462-2920.15935)  
882 [2920.15935](https://doi.org/10.1111/1462-2920.15935).
- 883 18. Das S, Bernasconi E, Koutsokera A, Wurlod D-A, Tripathi V, Bonilla-Rosso G,  
884 Aubert J-D, Derkenne M-F, Mercier L, Pattaroni C, Rapin A, von Garnier C, Marsland  
885 BJ, Engel P, Nicod LP. 2021. A prevalent and culturable microbiota links ecological  
886 balance to clinical stability of the human lung after transplantation. *Nat Commun* 12.
- 887 19. Auch AF, von Jan M, Klenk HP, Göker M. 2010. Digital DNA-DNA hybridization for  
888 microbial species delineation by means of genome-to-genome sequence comparison.  
889 *Stand Genomic Sci* 2:117.
- 890 20. Meier-Kolthoff JP, Göker M. 2019. TYGS is an automated high-throughput platform  
891 for state-of-the-art genome-based taxonomy. *Nat Commun* 10.
- 892 21. Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S. 2018. High  
893 throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries.  
894 *Nat Commun* 9:5114.
- 895 22. Beighton D, Hardie JM, Whitley RA. 1991. A scheme for the identification of viridans  
896 streptococci. *J Med Microbiol* 35:367–372.
- 897 23. Kawamura Y, Hou XG, Sultana F, Miura H, Ezaki T. 1995. Determination of 16S  
898 rRNA sequences of *Streptococcus mitis* and *Streptococcus gordonii* and phylogenetic  
899 relationships among members of the genus *Streptococcus*. *Int J Syst Bacteriol* 45:406–  
900 408.
- 901 24. Delorme C, Abraham AL, Renault P, Guédon E. 2015. Genomics of *Streptococcus*  
902 *salivarius*, a major human commensal. *Infection, Genetics and Evolution* 33:381–392.
- 903 25. Emms DM, Kelly S. 2019. OrthoFinder: Phylogenetic orthology inference for  
904 comparative genomics. *Genome Biol* 20:1–14.
- 905 26. Xu L, Dong Z, Fang L, Luo Y, Wei Z, Guo H, Zhang G, Gu YQ, Coleman-Derr D,  
906 Xia Q, Wang Y. 2019. OrthoVenn2: a web server for whole-genome comparison and  
907 annotation of orthologous clusters across multiple species. *Nucleic Acids Res* 47:W52.

- 908 27. Cantalapiedra CP, Hernández-Plaza A, Letunic I, Bork P, Huerta-Cepas J. 2021.  
909 eggNOG-mapper v2: Functional Annotation, Orthology Assignments, and Domain  
910 Prediction at the Metagenomic Scale. *Mol Biol Evol* 38:5825–5829.
- 911 28. Sarton-Lohéac G, Nunes da Silva CG, Mazel F, Baud G, de Bakker V, Das S, El  
912 Chazli Y, Ellegaard K, Garcia-Garcera M, Glover N, Liberti J, Nacif Marçal L, Prasad  
913 A, Somerville V, Bonilla-Rosso G, Engel P. 2023. Deep Divergence and Genomic  
914 Diversification of Gut Symbionts of Neotropical Stingless Bees. *mBio* 14:e0353822.
- 915 29. Price MN, Deutschbauer AM, Arkin AP. 2022. Filling gaps in bacterial catabolic  
916 pathways with computation and high-throughput genetics. *PLoS Genet* 18:e1010156.
- 917 30. Price MN, Deutschbauer AM, Arkin AP. 2020. GapMind: Automated Annotation of  
918 Amino Acid Biosynthesis. *mSystems* 5.
- 919 31. Zheng J, Ge Q, Yan Y, Zhang X, Huang L, Yin Y. 2023. dbCAN3: automated  
920 carbohydrate-active enzyme and substrate annotation. *Nucleic Acids Res* 51:W115–  
921 W121.
- 922 32. Abby SS, Néron B, Ménager H, Touchon M, Rocha EPC. 2014. MacSyFinder: A  
923 Program to Mine Genomes for Molecular Systems with an Application to CRISPR-  
924 Cas Systems. *PLoS One* 9:e110726.
- 925 33. Cantarel BI, Coutinho PM, Rancurel C, Bernard T, Lombard V, Henrissat B. 2009.  
926 The Carbohydrate-Active EnZymes database (CAZy): An expert resource for  
927 glycogenomics. *Nucleic Acids Res* 37.
- 928 34. Peterson SN, Sung CK, Cline R, Desai B V., Snesrud EC, Luo P, Walling J, Li H,  
929 Mintz M, Tsegaye G, Burr PC, Do Y, Ahn S, Gilbert J, Fleischmann RD, Morrison  
930 DA. 2004. Identification of competence pheromone responsive genes in *Streptococcus*  
931 *pneumoniae* by use of DNA microarrays. *Mol Microbiol* 51:1051–1070.
- 932 35. Whatmore AM, Barcus VA, Dowson CG. 1999. Genetic Diversity of the Streptococcal  
933 Competence (*com*) Gene Locus. *J Bacteriol* 181:3144–3154.
- 934 36. Luo P, Morrison DA. 2003. Transient association of an alternative sigma factor,  
935 ComX, with RNA polymerase during the period of competence for genetic  
936 transformation in *Streptococcus pneumoniae*. *J Bacteriol* 185:349–358.
- 937 37. Provvedi R, Dubnau D. 1999. ComEA is a DNA receptor for transformation of  
938 competent *Bacillus subtilis*. *Mol Microbiol* 31:271–280.
- 939 38. Diallo A, Foster HR, Gromek KA, Perry TN, Dujeancourt A, Krasteva P V., Gubellini  
940 F, Falbel TG, Burton BM, Fronzes R. 2017. Bacterial transformation: ComFA is a  
941 DNA-dependent ATPase that forms complexes with ComFC and DprA. *Mol*  
942 *Microbiol* 105:741–754.
- 943 39. Merritt J, Qi F, Shi W. 2005. A unique nine-gene *comY* operon in *Streptococcus*  
944 *mutans*. *Microbiology (Reading)* 151:157–166.
- 945 40. Haijema BJ, Hahn J, Haynes J, Dubnau D. 2001. A ComGA-dependent checkpoint  
946 limits growth during the escape from competence. *Mol Microbiol* 40:52–64.
- 947 41. Guglielmini J, Néron B, Abby SS, Garcillán-Barcia MP, La Cruz DF, Rocha EPC.  
948 2014. Key components of the eight classes of type IV secretion systems involved in  
949 bacterial conjugation or protein secretion. *Nucleic Acids Res* 42:5715.
- 950 42. Atmakuri K, Cascales E, Christie PJ. 2004. Energetic components VirD4, VirB11 and  
951 VirB4 mediate early DNA transfer reactions required for bacterial type IV secretion.  
952 *Mol Microbiol* 54:1199–1211.
- 953 43. Kim TJ, Bose N, Taylor RK. 2003. Secretion of a soluble colonization factor by the  
954 TCP type 4 pilus biogenesis pathway in *Vibrio cholerae*. *Mol Microbiol* 49:81–92.
- 955 44. Waksman G. 2019. From conjugation to T4S systems in Gram-negative bacteria: a  
956 mechanistic biology perspective. *EMBO Rep* 20:e47012.

- 957 45. Seemann T. GitHub - tseemann/abricate: :mag\_right: Mass screening of contigs for  
958 antimicrobial and virulence genes. <https://github.com/tseemann/abricate>. Retrieved 10  
959 November 2023.
- 960 46. Doster E, Lakin SM, Dean CJ, Wolfe C, Young JG, Boucher C, Belk KE, Noyes NR,  
961 Morley PS. 2020. MEGARes 2.0: a database for classification of antimicrobial drug,  
962 biocide and metal resistance determinants in metagenomic sequence data. *Nucleic  
963 Acids Res* 48:D561–D569.
- 964 47. Daly MM, Doktor S, Flamm R, Shortridge D. 2004. Characterization and Prevalence  
965 of MefA, MefE, and the Associated *msr* (D) Gene in *Streptococcus pneumoniae*  
966 Clinical Isolates. *J Clin Microbiol* 42:3570–3574.
- 967 48. Clancy J, Petitpas J, Dib-Hajj F, Yuan W, Cronan M, Kamath A V., Bergeron J,  
968 Retsema JA. 1996. Molecular cloning and functional analysis of a novel macrolide-  
969 resistance determinant, mefA, from *Streptococcus pyogenes*. *Mol Microbiol* 22:867–  
970 879.
- 971 49. Leclercq R, Courvalin P. 1991. Bacterial resistance to macrolide, lincosamide, and  
972 streptogramin antibiotics by target modification. *Antimicrob Agents Chemother*  
973 35:1267–1272.
- 974 50. Warburton PJ, Ciric L, Lerner A, Seville LA, Roberts AP, Mullany P, Allan E. 2013.  
975 TetAB(46), a predicted heterodimeric ABC transporter conferring tetracycline  
976 resistance in *Streptococcus australis* isolated from the oral cavity. *Journal of  
977 Antimicrobial Chemotherapy* 68:17.
- 978 51. Chen L, Yang J, Yu J, Yao Z, Sun L, Shen Y, Jin Q. 2005. VFDB: a reference  
979 database for bacterial virulence factors. *Nucleic Acids Res* 33.
- 980 52. Hu DK, Wang DG, Liu Y, Liu CB, Yu LH, Qu Y, Luo XH, Yang JH, Yu J, Zhang J,  
981 Li XY. 2013. Roles of virulence genes (PsaA and CpsA) on the invasion of  
982 *Streptococcus pneumoniae* into blood system. *Eur J Med Res* 18:1–6.
- 983 53. Hu Y, Park N, Seo KS, Park JY, Somarathne RP, Olivier AK, Fitzkee NC, Thornton  
984 JA. 2021. Pneumococcal surface adhesion A protein (PsaA) interacts with human  
985 Annexin A2 on airway epithelial cells. *Virulence* 12:1841–1854.
- 986 54. Carvalho MDGS, Tondella ML, McCaustland K, Weidlich L, McGee L, Mayer LW,  
987 Steigerwalt A, Whaley M, Facklam RR, Fields B, Carlone G, Ades EW, Dagan R,  
988 Sampson JS. 2007. Evaluation and Improvement of Real-Time PCR Assays Targeting  
989 *lytA*, *ply*, and *psaA* Genes for Detection of Pneumococcal DNA. *J Clin Microbiol*  
990 45:2460–2466.
- 991 55. Skov Sørensen UB, Yao K, Yang Y, Tettelin H, Kilian M. 2016. Capsular  
992 polysaccharide expression in commensal *Streptococcus* species: Genetic and antigenic  
993 similarities to *Streptococcus pneumoniae*. *mBio* 7.
- 994 56. Agapov VS, Smirenskaia T V, Komnova ZD. 1987. [Clinico-morphological  
995 characteristics of periradicular cysts bordering on the maxillary sinus]. *Stomatologia*  
996 (Mosk) 66:11–3.
- 997 57. Guerin ME, Kordulakova J, Schaeffer F, Svetlikova Z, Buschiazzo A, Giganti D,  
998 Gicquel B, Mikusova K, Jackson M, Alzari PM. 2007. Molecular recognition and  
999 interfacial catalysis by the essential phosphatidylinositol mannosyltransferase PimA  
1000 from mycobacteria. *J Biol Chem* 282:20705–20714.
- 1001 58. Slager J, Aprianto R, Veening JW. 2018. Deep genome annotation of the opportunistic  
1002 human pathogen *Streptococcus pneumoniae* D39. *Nucleic Acids Res* 46:9971–9989.
- 1003 59. Mohanty P, Patel A, Bhardwaj AK. 2012. Role of H- and D- MATE-type transporters  
1004 from multidrug resistant clinical isolates of *Vibrio fluvialis* in conferring  
1005 fluoroquinolone resistance. *PLoS One* 7.

- 1006 60. Smith HE, Damman M, van der Velde J, Wagenaar F, Wisselink HJ, Stockhofe-  
1007 Zurwieden N, Smits MA. 1999. Identification and Characterization of the *cps* Locus of  
1008 *Streptococcus suis* Serotype 2: the Capsule Protects against Phagocytosis and Is an  
1009 Important Virulence Factor. *Infect Immun* 67:1750–1756.
- 1010 61. Escapa IF, Chen T, Huang Y, Gajare P, Dewhirst FE, Lemon KP. 2018. New Insights  
1011 into Human Nostril Microbiome from the Expanded Human Oral Microbiome  
1012 Database (eHOMD): a Resource for the Microbiome of the Human Aerodigestive  
1013 Tract. *mSystems* 3.
- 1014 62. Segal LN, Alekseyenko A V., Clemente JC, Kulkarni R, Wu B, Chen H, Berger KI,  
1015 Goldring RM, Rom WN, Blaser MJ, Weiden MD. 2013. Enrichment of lung  
1016 microbiome with supraglottic taxa is associated with increased pulmonary  
1017 inflammation. *Microbiome* 1:19.
- 1018 63. Erb-Downward JR, Thompson DL, Han MK, Freeman CM, McCloskey L, Schmidt  
1019 LA, Young VB, Toews GB, Curtis JL, Sundaram B, Martinez FJ, Huffnagle GB. 2011.  
1020 Analysis of the lung microbiome in the “healthy” smoker and in COPD. *PLoS One*  
1021 <https://doi.org/10.1371/journal.pone.0016384>.
- 1022 64. Charlson ES, Diamond JM, Bittinger K, Fitzgerald AS, Yadav A, Haas AR, Bushman  
1023 FD, Collman RG. 2012. Lung-enriched Organisms and Aberrant Bacterial and Fungal  
1024 Respiratory Microbiota after Lung Transplant. *Am J Respir Crit Care Med* 186:536–  
1025 545.
- 1026 65. Borewicz K, Pragman AA, Kim HB, Hertz M, Wendt C, Isaacson RE. 2013.  
1027 Longitudinal analysis of the lung microbiome in lung transplantation. *FEMS Microbiol*  
1028 *Lett* <https://doi.org/10.1111/1574-6968.12053>.
- 1029 66. Sharma NS, Vestal G, Wille K, Patel KN, Cheng F, Tipparaju S, Tousif S, Banday  
1030 MM, Xu X, Wilson L, Nair VS, Morrow C, Hayes D, Seyfang A, Barnes S, Deshane  
1031 JS, Gaggari A. 2020. Differences in airway microbiome and metabolome of single lung  
1032 transplant recipients. *Respir Res* 21:1–12.
- 1033 67. Doern CD, Carey-Ann BD. 2010. It’s Not Easy Being Green: the Viridans Group  
1034 Streptococci, with a Focus on Pediatric Clinical Manifestations. *J Clin Microbiol*  
1035 48:3829–3835.
- 1036 68. Ruoff KL. 1991. Nutritionally Variant Streptococci. *Clin Microbiol Rev* 4:184–190.
- 1037 69. Kim JN, Ahn SJ, Burne RA. 2015. Genetics and Physiology of Acetate Metabolism by  
1038 the Pta-Ack Pathway of *Streptococcus mutans*. *Appl Environ Microbiol* 81:5015.
- 1039 70. Garnett JP, Baker EH, Baines DL. 2012. Sweet talk: insights into the nature and  
1040 importance of glucose transport in lung epithelium. *European Respiratory Journal*  
1041 40:1269–1276.
- 1042 71. Kwong WK, Zheng H, Moran NA. 2018. Erratum to: Convergent evolution of a  
1043 modified, acetate-driven TCA cycle in bacteria (*Nature Microbiology*, (2017), 2,  
1044 (17067), 10.1038/nmicrobiol.2017.67). *Nat Microbiol* 3:960.
- 1045 72. Sawers RG, Clark DP. 2004. Fermentative Pyruvate and Acetyl-Coenzyme A  
1046 Metabolism. *EcoSal Plus* 1.
- 1047 73. Tagaino R, Washio J, Abiko Y, Tanda N, Sasaki K, Takahashi N. 2019. Metabolic  
1048 property of acetaldehyde production from ethanol and glucose by oral *Streptococcus*  
1049 and *Neisseria*. *Scientific Reports* 2019 9:1 9:1–8.
- 1050 74. Deleu S, Machiels K, Raes J, Verbeke K, Vermeire S. 2021. Short chain fatty acids  
1051 and its producing organisms: An overlooked therapy for IBD? *EBioMedicine*  
1052 66:103293.
- 1053 75. Antunes KH, Fachi JL, de Paula R, da Silva EF, Pral LP, dos Santos AÁ, Dias GBM,  
1054 Vargas JE, Puga R, Mayer FQ, Maito F, Zárata-Bladés CR, Ajami NJ, Sant’Ana MR,  
1055 Candreva T, Rodrigues HG, Schmiele M, Silva Clerici MTP, Proença-Modena JL,

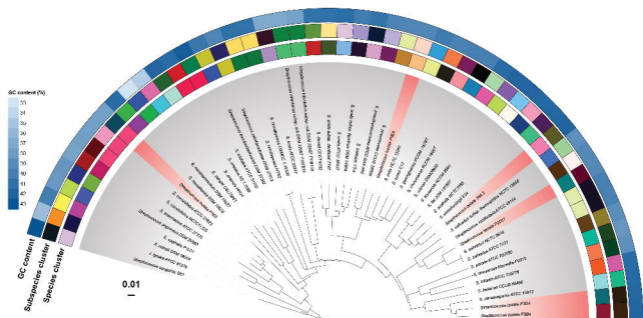
- 1056 Vieira AT, Mackay CR, Mansur D, Caballero MT, Marzec J, Li J, Wang X, Bell D,  
1057 Polack FP, Kleeberger SR, Stein RT, Vinolo MAR, de Souza APD. 2019. Microbiota-  
1058 derived acetate protects against respiratory syncytial virus infection through a GPR43-  
1059 type 1 interferon response. *Nat Commun* 10:3273.
- 1060 76. Machado MG, Patente TA, Rouillé Y, Heumel S, Melo EM, Deruyter L, Pourcet B,  
1061 Sencio V, Teixeira MM, Trottein F. 2022. Acetate Improves the Killing of  
1062 *Streptococcus pneumoniae* by Alveolar Macrophages via NLRP3 Inflammasome and  
1063 Glycolysis-HIF-1 $\alpha$  Axis. *Front Immunol* 13.
- 1064 77. Håvarstein LS, Coomaraswamy G, Morrison DA. 1995. An unmodified  
1065 heptadecapeptide pheromone induces competence for genetic transformation in  
1066 *Streptococcus pneumoniae*. *Proc Natl Acad Sci U S A* 92:11140.
- 1067 78. Haenni M, Lupo A, Madec J-Y. 2018. Antimicrobial Resistance in *Streptococcus* spp .  
1068 *Microbiol Spectr* 6.
- 1069 79. Gibson PS, Veening J-W. 2023. Gaps in the wall: understanding cell wall biology to  
1070 tackle amoxicillin resistance in *Streptococcus pneumoniae*. *Curr Opin Microbiol*  
1071 72:102261.
- 1072 80. Skov Sørensen UB, Yao K, Yang Y, Tettelin H, Kilian M. 2016. Capsular  
1073 Polysaccharide Expression in Commensal *Streptococcus* Species: Genetic and  
1074 Antigenic Similarities to *Streptococcus pneumoniae*. *mBio* 7.
- 1075 81. Chang B, Morita M, Nariai A, Kasahara K, Kakutani A, Ogawa M, Ohnishi M, Oishi  
1076 K. 2022. Invasive *Streptococcus oralis* Expressing Serotype 3 Pneumococcal Capsule,  
1077 Japan - Volume 28, Number 8—August 2022 - *Emerging Infectious Diseases journal* -  
1078 CDC. *Emerg Infect Dis* 28:1720–1722.
- 1079 82. s-andrews/FastQC: A quality control analysis tool for high throughput sequencing  
1080 data. <https://github.com/s-andrews/FastQC>. Retrieved 13 July 2023.
- 1081 83. Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina  
1082 sequence data. *Bioinformatics* 30:2114–2120.
- 1083 84. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM,  
1084 Nikolenko SI, Pham S, Prjibelski AD, Pyshkin A V., Sirotkin A V., Vyahhi N, Tesler  
1085 G, Alekseyev MA, Pevzner PA. 2012. SPAdes: a new genome assembly algorithm and  
1086 its applications to single-cell sequencing. *J Comput Biol* 19:455–477.
- 1087 85. Mikheenko A, Saveliev V, Gurevich A. 2016. MetaQUAST: evaluation of  
1088 metagenome assemblies. *Bioinformatics* 32:1088–1090.
- 1089 86. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. 2015. CheckM:  
1090 assessing the quality of microbial genomes recovered from isolates, single cells, and  
1091 metagenomes. *Genome Res* 25:1043–1055.
- 1092 87. Seemann T. 2014. Prokka: rapid prokaryotic genome annotation. *Bioinformatics*  
1093 30:2068–2069.
- 1094 88. Letunic I, Bork P. 2007. Interactive Tree Of Life (iTOL): an online tool for  
1095 phylogenetic tree display and annotation. *Bioinformatics* 23:127–128.
- 1096 89. Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S. 2018. High  
1097 throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries.  
1098 *Nature Communications* 2018 9:1 9:1–8.
- 1099 90. aniclustermmap · PyPI. <https://pypi.org/project/aniclustermmap/>. Retrieved 15 June 2023.
- 1100 91. Katoh K, Misawa K, Kuma KI, Miyata T. 2002. MAFFT: a novel method for rapid  
1101 multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res*  
1102 30:3059–3066.
- 1103 92. Price MN, Dehal PS, Arkin AP. 2009. FastTree: Computing Large Minimum  
1104 Evolution Trees with Profiles instead of a Distance Matrix. *Mol Biol Evol* 26:1641–  
1105 1650.

- 1106 93. Stamatakis A. 2006. RAxML-VI-HPC: maximum likelihood-based phylogenetic  
1107 analyses with thousands of taxa and mixed models. *Bioinformatics* 22:2688–2690.
- 1108 94. Mostacci N, Wüthrich TM, Siegwald L, Kieser S, Steinberg R, Sakwinska O, Latzin P,  
1109 Korten I, Hilty M. 2023. Informed interpretation of metagenomic data by StrainPhlAn  
1110 enables strain retention analyses of the upper airway microbiome. *mSystems*  
1111 8:e0072423.
- 1112 95. Neufeld F. 1902. Ueber die Agglutination der Pneumokokken und über die Theorieen  
1113 der Agglutination. *Zeitschrift für Hygiene und Infektionskrankheiten* 1902 40:1 40:54–  
1114 72.
- 1115 96. Agapov VS, Smirenskaia T V, Komnova ZD. 1987. [Clinico-morphological  
1116 characteristics of periradicular cysts bordering on the maxillary sinus]. *Stomatologia*  
1117 (Mosk) 66:11–3.
- 1118 97. Yin Y, Mao X, Yang J, Chen X, Mao F, Xu Y. 2012. dbCAN: a web resource for  
1119 automated carbohydrate-active enzyme annotation. *Nucleic Acids Res* 40:W445–  
1120 W451.
- 1121 98. Drula E, Garron ML, Dogan S, Lombard V, Henrissat B, Terrapon N. 2022. The  
1122 carbohydrate-active enzyme database: functions and literature. *Nucleic Acids Res*  
1123 50:D571–D577.
- 1124 99. Abby SS, Néron B, Ménager H, Touchon M, Rocha EPC. 2014. MacSyFinder: A  
1125 Program to Mine Genomes for Molecular Systems with an Application to CRISPR-  
1126 Cas Systems. *PLoS One* 9:e110726.
- 1127 100. Shimoyama Y. 2022. pyGenomeViz: A genome visualization python package for  
1128 comparative genomics. <https://github.com/moshi4/pyGenomeViz>. Retrieved 14  
1129 December 2023.
- 1130 101. Edgar RC. 2004. MUSCLE: a multiple sequence alignment method with reduced time  
1131 and space complexity. *BMC Bioinformatics* 5:113.
- 1132

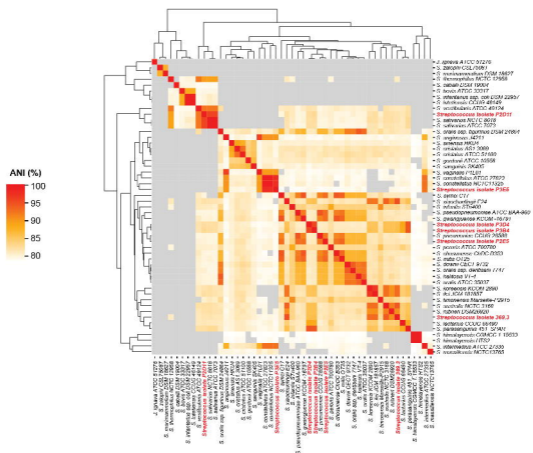


# Figure 1

## A Whole genome BLAST Distance Phylogeny

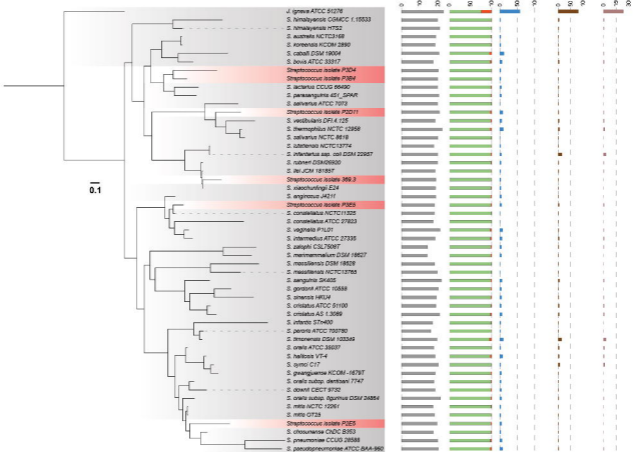


## B Pairwise Average Nucleotide Identity (ANI)

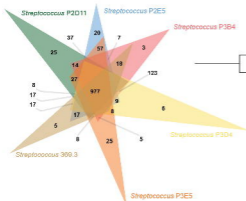


**Figure 2**

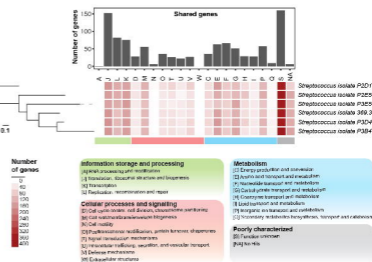
**A** Core genome phylogeny and orthologous genes



**B** Shared genes in lung streptococci

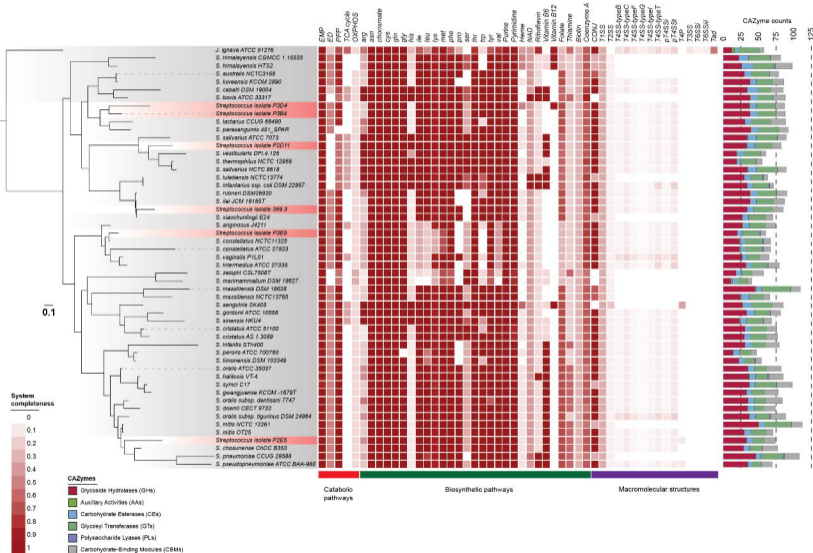


**C** Functional classification of COGs



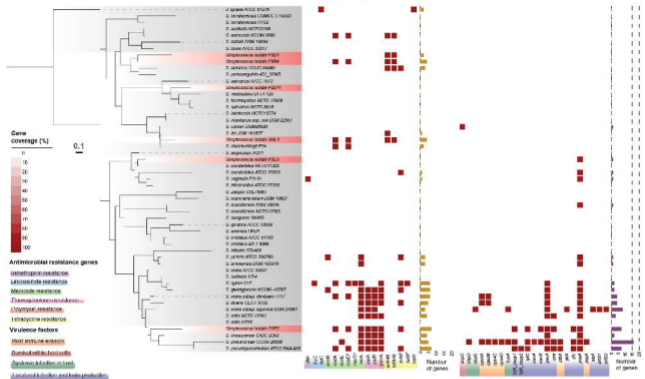
# Figure 3

## Metabolic pathways and macromolecular structures in lung streptococci

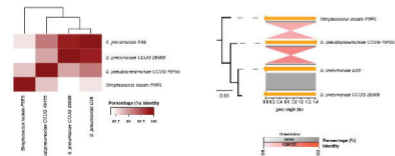


# Figure 4

## A Antimicrobial resistance and virulence factors in human lung streptococci

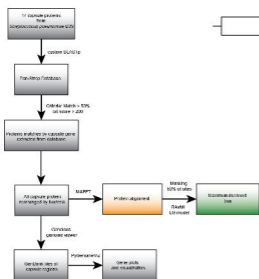


## B Pneumolysin gene in the novel lung *Streptococcus* P2E5

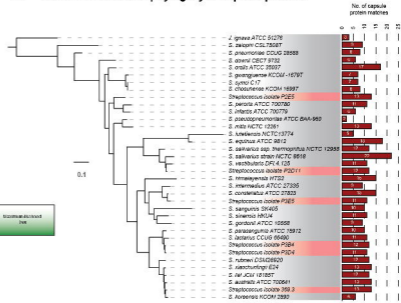


# Figure 9

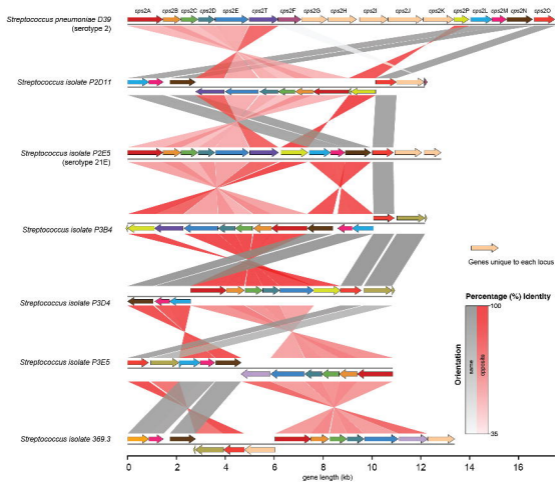
## A Capsule analysis pipeline



## B Maximum-likelihood phylogeny of capsule proteins



## C Genomic arrangement and synteny of capsule genes



**Figure 6**

Maximum likelihood tree showing evolutionary relationship based on single-copy core proteins

